

Judgment sampling

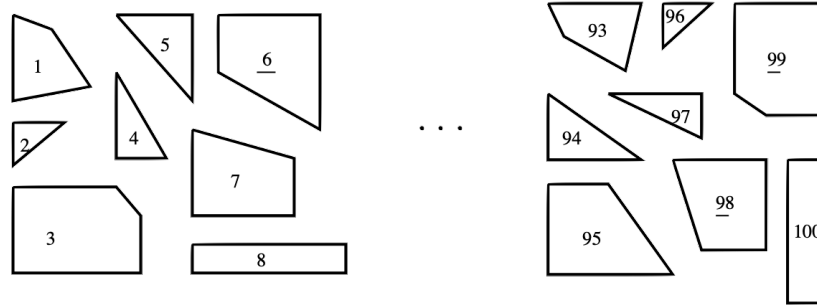
```
## Set this up for your own directory
imageDirectory <- "MyAssignmentDirectory/img" # e.g. in current "./img"
dataDirectory <- "MyAssignmentDirectory/data" # e.g. in current "./data"
path_concat <- function(path1, path2, sep="/") paste(path1, path2, sep = sep)
```

27 marks

A number of graduate data science students were presented with the following competition.

Consider a study population \mathcal{P}_{Study} consisting of $N = 100$ blocks labelled $u = 1, 2, 3, \dots, 100$.

The blocks are of uniform thickness and density (all blocks were cut from the same opaque plastic sheet of about 5mm thickness), but have different shapes as such as shown below:



Suppose also that $\mathcal{P}_{Target} = \mathcal{P}_{Study}$ and that the population attribute of interest

$$a(\mathcal{P}_{Target}) = \frac{1}{N} \sum_{u \in \mathcal{P}_{Target}} weight(u)$$

that is the average weight of all $N = 100$ blocks in the population.

We want a sample $\mathcal{S} \subset \mathcal{P}_{Study}$ of $n = 10$ blocks selected from the 100, whose average weight is (nearly) the same as the average weight of all 100.

That is, we would like a sample with zero (or at least small in absolute value) **sample error** $a(\mathcal{S}) - a(\mathcal{P}_{Study})$.

The `blocks` data can be loaded from the assignment data directory as follows:

```
load(path_concat(dataDirectory, "blocks.rda"))
head(blocks, n = 3)
```

```
##   id weight perimeter group
## 1  1     55         32     B
## 2  2     35         27     B
## 3  3     35         25     A
```

Having been presented with all 100 blocks and asked to **judge** which 10 blocks have an average weight nearest the average weight of all 100 blocks, each student would have come up with their own sampling plan based on their judgment. This type of sampling is called **judgment sampling**.

The id numbers of the students and the blocks they selected are recorded in another file, `judgmentSamples.csv`. These can be loaded from the assignment data directory as follows:

```
students <- read.csv(path_concat(dataDirectory, "judgmentSamples.csv"))
head(students, n = 3)
```

```
## studentID first second third fourth fifth sixth seventh eighth ninth tenth
## 1      5086    12     18    17     11    15    20     14     13     16     18
## 2      3848    34     35    70     56    32    14      5     88     81     73
## 3      6656    14     34    41     29    32    55     74     40     16     70
```

The variates of `student` identify the student and the id numbers of the blocks they selected, in the order they recorded them.

- a. (4 marks) Draw a histogram of all of the block weights selected by the students. If any block was selected by more than one student, include its weight as often as it was selected. That is, there will be a total of 330 weights used to construct the histogram.
 - Make sure the histogram is suitably labelled
 - Add a vertical dashed red line at the at the average of all 100 weights in the entire population of 100 blocks (i.e. not just those selected by students).

Show your code.

- b. (5 marks) For each student, calculate the sample average weight of the blocks they selected. Create a data frame called `judgmentErrors` of the student ids and their sample errors. Print out the ids and sample errors for both the top five and the bottom five students in increasing order of their *absolute* sample error.

Show your code.

- c. (3 marks) Estimate the sampling bias and the sampling standard deviation for judgment sampling on this data. Show your code.

- d. (3 marks) Provide a (suitably labelled) histogram of the sample errors. Add a vertical red dashed line at 0.

- e. (3 marks) Calculate the sample standard deviation of the weights selected for each of the judgment samples. Draw a histogram of these standard deviations (suitably labelled). Draw a vertical dashed red line at the average of these standard deviations.

Show your code.

- f. (6 marks) Identify which student had the smallest sample standard deviation **and** which student had the largest sample standard deviation. Report their standard deviations. Draw histograms (suitably labelled **and** having the same `xlim = extendrange(blocks$weight)`) of the weights of the blocks selected by each of these students. Add a vertical dashed red line to each histogram at the average of all 100 block weights in the population. What do you conclude about the sampling plan of each of these students?

Show your code.

- g. (3 marks) Comment on the quality of this judgment sampling plan, making reference to any of the results calculated above.